

Accelerating Unstructured Mesh Computations using FPGAs

Kyrylo Tkachov

Supervisor: Professor Paul Kelly

May 27, 2012

Abstract

In this report we present a methodology for accelerating computations performed on unstructured meshes in the course of a finite volume approach. We use Field Programmable Gate Arrays, or FPGAs, to create an application-specific datapath that will be used as a co-processor to perform the bulk of the floating point operations required by the application. In particular, we focus on dealing with irregular memory access patterns that are a consequence of using an unstructured mesh in such a way as to facilitate a streaming model of computation. We describe the partitioning of the mesh and the techniques used to exchange information between neighbouring partitions, using so-called halos. We provide an implementation of a concrete 2D finite volume application and consider the extension to 3D and more complex computations. We evaluate our results by comparing the speedup achieved with analogous GPGPU and multi-core processor implementations.

Acknowledgements

I would like to thank Professor Paul Kelly for giving me so much of his time and ideas and making me aware of scope of the project and the intricacies involved. I would also like to thank Dr. Carlo Bertolli for providing practical advice and explaining the labyrinth that is heterogenous computing. Special thanks go to the team at Maxeler Technologies for helping me out with the details of FPGA-based acceleration and providing support for their excellent toolchain I extend my gratitude to Dr. Tony Field, my personal tutor, who supported me throughout my years at Imperial College and guided so much of my academic development, as well as being the second supervisor on this project.

I would like to thank my mother and grandfather for supporting me through university, both materially and psychologically. Last but not least, I would like to thank my coursemates and friends all over the world, with whom I've had many thought-provoking discussions on every subject imaginable and who always kept me motivated, even when I doubted myself.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 3 |
| 1.1 | The domain | 3 |
| 1.2 | The Airfoil program | 4 |
| 1.3 | FPGAs, streaming and acceleration | 4 |
| 1.4 | Contributions | 6 |
| 2 | Background | 8 |
| 2.1 | Unstructured meshes and their representation | 8 |
| 2.2 | Airfoil | 10 |
| 2.2.1 | Computational kernels and data sets | 12 |
| 2.2.2 | Indirection maps | 14 |
| 2.3 | Hardware platform, Maxeler toolchain and the streaming model of computation | 17 |
| 2.3.1 | MaxCompiler example | 18 |
| 2.3.2 | Hardware | 24 |
| 2.3.3 | Mesh Partitioning and halos | 24 |
| 2.4 | Previous work | 25 |
| 3 | Details and implementation | 28 |
| 3.1 | Plan | 28 |
| 4 | Appendix | 32 |
| 4.1 | Airfoil Kernel definitions in C | 32 |

Chapter 1

Introduction

This project presents a methodology for accelerating computations performed on unstructured meshes in the context of Computational Fluid Dynamics (CFD). We use Field Programmable Gate Arrays, or FPGAs, to construct a high-throughput streaming pipeline which is kept filled thanks to an appropriate data layout and partitioning scheme for the mesh. We explore the rearrangement and grouping schemes used to achieve locality of the data points. A formal performance model is constructed to predict the performance characteristics of our architecture and hence justify the design choices made. Appropriate evaluation tests are performed to evaluate the results on a sample CFD application, achieving speedup comparable with state of the art GPGPU and multi-processor solutions. In this section we present a general overview of the problem domain, the hardware platform and the contributions of this project.

1.1 The domain

Computational Fluid Dynamics, or CFD, is a branch of physics focused on numerical algorithms that simulate the movement of fluids and gases and their interactions with surfaces. These simulations are widely used by engineers to design structures and equipment that interact with fluid substances, for example airplane wings and turbines, water and oil pipelines etc.

The required calculations are usually expressed as systems of partial differential equations, the Navier-Stokes equations or the Euler equations, which are discretized using any of a number of techniques. The technique used by our sample application, Airfoil, is the finite volume method that

calculates values at discrete places in a mesh and relies on the observation that the fluxes entering a volume are equal to the fluxes leaving it. This project is not concerned with the exact mathematical formulation of these techniques, but they provide a feel for the origins of the problem domain.

1.2 The Airfoil program

The sample program we examine is called Airfoil, a 2D unstructured mesh finite volume simulation of fluid motion around an airplane wing (which has the shape of an airfoil). Airfoil was written as a representative of the class of programs that are tackled by OP2, a framework partially developed and maintained by the Software Performance Optimisation group at Imperial College to abstract the acceleration of unstructured mesh computations on a wide variety of hardware backends.

Airfoil defines an unstructured mesh through sets of nodes, edges and cells and associating them through mappings. Airfoil is written in the C language and these sets are represented at the lowest level as C-arrays. Then data is associated with these sets, such as node coordinates, temperature, pressure etc. The mesh solution is then expressed as the conceptually parallel application of computational kernels on the data associated with each element of a particular set (nodes, edges, cells). These kernels are usually floating point- intensive operations and update the datasets. The procedure is repeated through multiple iterations as desired until a steady-state solution is reached. A more detailed discussion of the unstructured mesh is presented in the Background section of this report.

1.3 FPGAs, streaming and acceleration

In this project we explore the acceleration possibilities of problems in the described domain by using Field Programmable Gate Arrays, or FPGAs. FPGAs are integrated circuits that can be reconfigured on the fly to implement in hardware any logic design. Thanks to this property they provide the development flexibility of software with the benefits of an explicit custom hardware datapath. At a high level, FPGAs can be viewed as a two-dimensional grid of logic elements that can be interconnected in any desirable way.

The FPGA acceleration approach we look at is the streaming model of computation. In a streaming approach we create a dataflow graph out of simple computational nodes that perform a specific operation on pieces of

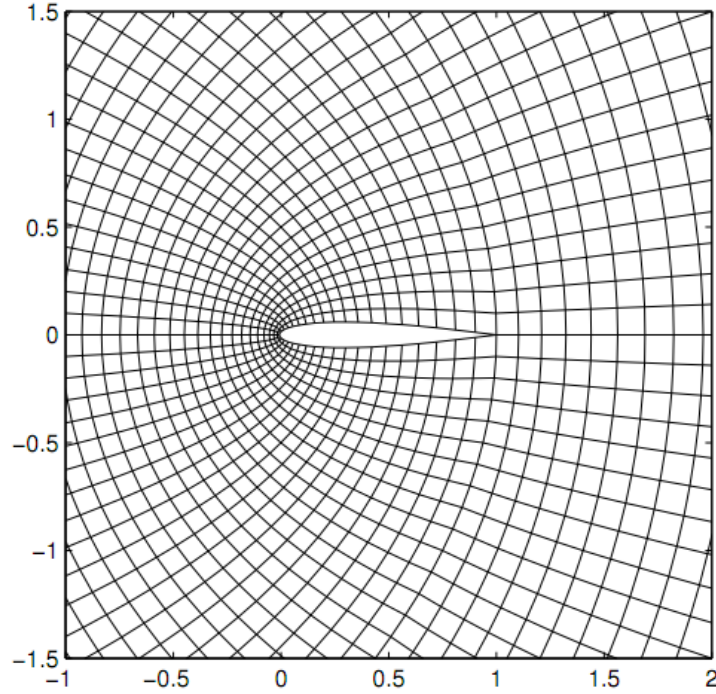


Figure 1.1: Visualisation of a reduced version of the Airfoil mesh

data pushed in and out of them. Connecting these nodes together creates a pipeline through which one can stream an array of data and get one output per cycle thus achieving high throughput. A simple dataflow graph can be seen in figure 1.2. FPGAs are usually programmed using a low level hardware description language like VHDL or Verilog, however many tools have been designed that allow a developer to specify high-level designs. We use MaxCompiler, a compiler that lets us specify the computational graph through a high-level Java API, so we focus on the functional aspects of our design and the tool generates a hardware implementation of it. We use this approach to implement a datapath the kernel described in Airfoil and we then look at approaches to utilise the streaming bandwidth. The FPGAs we consider have a large DRAM storage area attached (24GB) to them that can be used to store the mesh and utilising the bandwidth of that DRAM fully is key to achieving maximum performance.

During the course of our work it emerges that in order to stream data to and from the accelerator continuously, we need to enforce some spatial locality in the mesh data, thus requiring us to reorder the data and or-

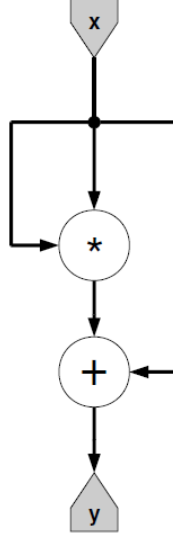


Figure 1.2: A simple dataflow graph that implements the function $y(x) = x^2 + x$.

ganise it into partitions that will be stored in the kernel internally and will need to exchange data with neighbouring partitions through a halo exchange mechanism. This opens a whole new space of decisions that we must make pertaining to the storage layout and streaming responsibilities of the DRAM and the host machine. We present the mesh partitioning schemes that are used to maximise DRAM bandwidth utilisation and maximise pipelining.

We present a performance model that will be used to describe the theoretical performance increase of the system in terms of various parameters like DRAM utilisation, clock frequency etc. Finally we evaluate the performance of our implementation of Airfoil against existing GPGPU and multi-processors cluster implementations.

1.4 Contributions

- We present a methodology for accelerating unstructured mesh computations using deeply pipelined streaming FPGA designs.
- We investigate memory layout issues that arise from efforts to maximise the spatial locality of the mesh.

- We provide a hardware accelerated version of part of the Airfoil program using the methodologies described in this report.
- We provide a predictive performance model that is used to justify our design decisions and provide a formal expression of the potential speedup.
- We investigate the potential for generalisation of the problem and the acceleration of more complex industry-grade unstructured mesh simulations.

Chapter 2

Background

This section provides more detail on the sample application. An overview of the Maxeler toolchain is given, which is used to implement the streaming solution we develop. The streaming model of computation is presented in the context of MaxCompiler by walking through steps to build a simple MaxCompiler application. Previous work in this area is presented and summarised in order to provide a context for the contributions of our work.

2.1 Unstructured meshes and their representation

The spatial domain of the problem can be discretised into either a structured or an unstructured mesh. A structured mesh has the advantage of having a highly regular structure and thus a highly predictable access pattern. If, however, one needs a more detailed solution around a particular area, the mesh would have to be fine-grained across the whole domain, thus increasing the number of cells, nodes and edges by a large factor even in areas that are not of such great interest. This is where unstructured meshes come in. They explicitly describe the connectivity between the elements and can thus be refined and coarsened around particular areas of interest. This provides much greater flexibility at the expense of losing the regularity of the mesh, forcing us to store the connectivity information that defines its topology. It is useful to have an intimate understanding of the representation of unstructured meshes in order to understand the techniques discussed further on. A graphical example is shown in figure 2.1

In our sample application the mesh distinguishes three main elements: nodes, cells and edges. We have to represent the connectivity information between them. This is done through *indirection maps* that store, for exam-

ple, the nodes that an edge connects or the nodes that a cell contains. In the application we explore the cells always have four nodes and the edges always connect two nodes and have two cells adjacent. In the more general case of variable-dimension cells (quadrilaterals, triangles, hexagons all mixed together) we would need an additional array storing the indices into the indirection maps and the sizes of the elements. But we do not consider such meshes here.

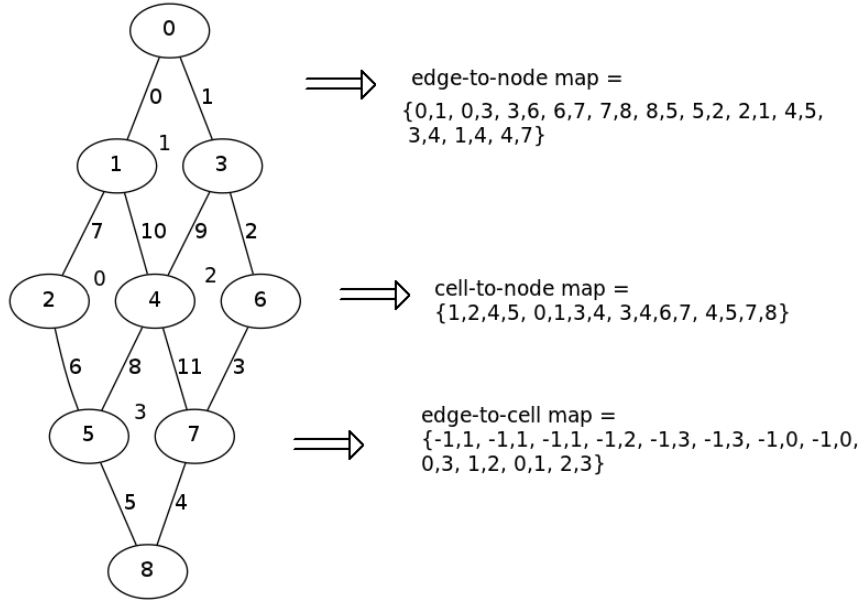


Figure 2.1: An example mesh and its representation using indirection arrays. The cell numbers are shown inside the quadrilaterals formed by the nodes (circles) and edges (edges connecting the nodes). Together with the indirection map, we also store an integer $dim \in \mathbb{N}$ which specifies the dimension of the mapping. Thus, the data associated with element i are stored in the range $[i * dim, \dots, i * (dim + 1) - 1]$ of the relevant indirection map (in the example: the nodes associated with edge 3 are stored at indices $3 * 2 = 6$ and $3 * 2 + 1 = 7$). Note: in the edge-to-cell map -1 represents a boundary cell that may be handled in a special way by a computational kernel.

The above method deals with the connectivity information amongst the different elements of the mesh. The data on which we perform the actual arithmetic calculations is stored in arrays indexed by element number. Such an approach is presented in figure 2.2.

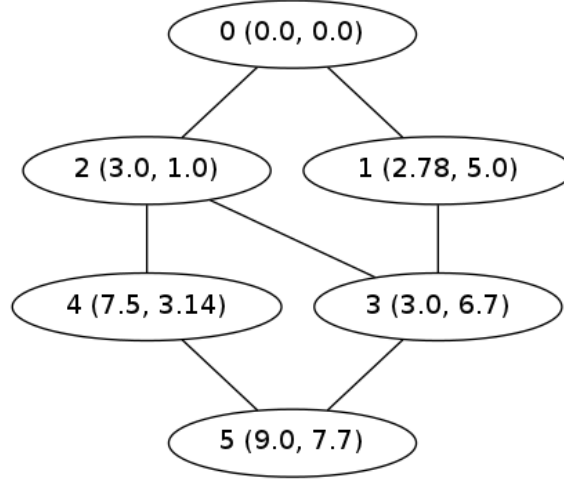


Figure 2.2: An example mesh with coordinate data associated with each node $((x, y)$ from $node_id$ (x, y)). The coordinate data will be represented as an array of floating point numbers $x = \{0.0, 0.0, 2.78, 5.0, 3.0, 1.0, 3.0, 6.7, 7.5, 3.14, 9.0, 7.7\}$. Again we also record the dimension of the data (in this case $dim = 2$) in order to access the data set associated with each element. In this example, the coordinate data for node 4 is stored at indices $4 * 2 = 8$ and $4 * 2 + 1 = 9$ of the array x .

2.2 Airfoil

Airfoil was written as a representative of the class of problems we are interested in. It was initially designed as a non-trivial example of the issues tackled by the OP2 framework. Although we are not directly dealing with OP2 in this project, an overview of Airfoil within this context is provided by MB Giles et al [1] because it discusses the acceleration issues arising from the memory access pattern.

The computational work in Airfoil is performed by 5 loops that work one after the other and operate on the nodes, cells and edges of the mesh. They work by applying a kernel on the data item referenced by the node, edge or cell. Conceptually, the application of a kernel to a data item is independent of the application to any other item in the same set, and can therefore be executed in parallel. The complexity comes from reduce operations, where some edges or cells update the same data item (associated with the same node). In these cases care must be taken to ensure the correct update of the data. For parallel architectures such as GPUs and multi-processor clusters

this issue can be resolved by enforcing an atomic commit scheme or by colouring the mesh partitions, so that no two partitions update the same data item simultaneously [1].

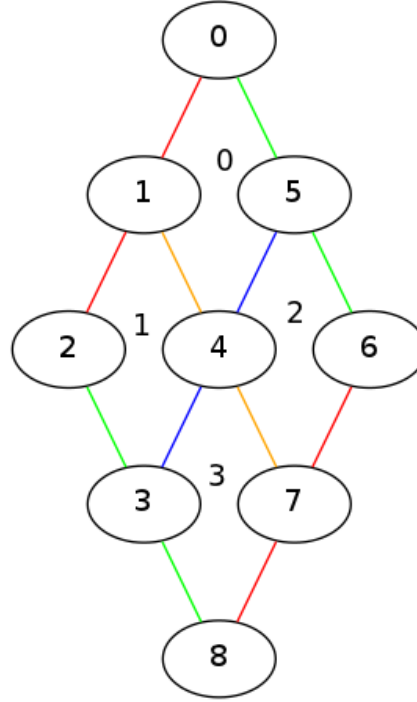


Figure 2.3: An example mesh, showing data dependencies between edges that affect cell data.

Consider figure 2.3. Take for example edges $\alpha = (1, 4)$ and $\beta = (4, 5)$. Say there is a data item x associated with every cell and the processing of an edge increments the data items associated with its two cells. α and β cannot execute in parallel because they are both associated with cell 4 and can therefore end up using out of date copies of the data associated with cell 4 by the following sequence of events: α reads initial x_0 , β reads x_0 , α computes $x_\alpha = x_0 + 1$, β computes $x_\beta = x_0 + 1$, α writes back x_α , β writes back x_β and the final value of x turns out to be $x_\beta = x_0 + 1$ instead of the desired $x_0 + 2$. Some implementations work around this issue by colouring the edges, such that no two edges of the same colour share a cell and can therefore be processed in parallel. Figure 2.3 shows such a colouring. Another option would be to introduce atomic operations and/or locking, but

that would approach severely limits parallelisation opportunities.

2.2.1 Computational kernels and data sets

Airfoil defines five computational kernels that iterate over the mesh, performing floating point calculations on the data sets defined over the elements of the mesh. We shall describe them by the elements they iterate over and by the elements they read and modify. As described above, we also define some data sets that are associated with the mesh elements. The datasets defined in Airfoil are shown in table 2.1.

| Data set name | Associated with | Type/Dimension | Physical meaning |
|---------------|-----------------|--|--|
| x | Nodes | $\mathbb{R} \times \mathbb{R}$ | Node coordinates |
| q | Cells | $\mathbb{R} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}$ | density, momentum, energy per unit volume |
| q_old | Cells | $\mathbb{R} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}$ | values of q from previous iteration |
| res | Cells | $\mathbb{R} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}$ | residual |
| adt | Cells | \mathbb{R} | Used for calculating area/timestep |
| bound | Edges | $\{0, 1\}$ | Specifies whether an edge is on the boundary of the mesh |

Table 2.1: Table showing the data sets and their types. In the actual implementation, we may choose to represent real numbers (\mathbb{R}) as standard or double precision floating point numbers or as fixed point numbers (discussed later). Elements of dimension larger than one will be represented as arrays. The physical meaning of these sets is not important, however Airfoil is generally interested in computing a steady-state solution for the q data set.

The kernels are presented in table 2.2 along with the datasets they require and modify. To show a more concrete example of what these kernels do, the `res_calc` kernel code in the C language is presented in Listing 1. The rest of the kernels are reproduced in the appendix

| Kernel Name | Iterates over | Reads | Writes |
|-------------|------------------|------------------|--------|
| save_soln | Cells | q | q_old |
| adt_calc | Cells | x, q | adt |
| res_calc | Edges | x, q, adt | res |
| bres_calc | (Boundary) Edges | x, q, adt, bound | res |
| update | Cells | q_old, adt, res | q, res |

Table 2.2: Table showing the kernels defined in airfoil and their data requirements.

```

1  void res_calc(float *x1, float *x2, float *q1, float *q2,
2              float *adt1, float *adt2, float *res1, float *res2) {
3      float dx,dy,mu, ri, p1,vol1, p2,vol2, f;
4      dx = x1[0] - x2[0];
5      dy = x1[1] - x2[1];
6      ri = 1.0f/q1[0];
7      p1 = gm1*(q1[3]-0.5f*ri*(q1[1]*q1[1]+q1[2]*q1[2]));
8      vol1 = ri*(q1[1]*dy - q1[2]*dx);
9      ri = 1.0f/q2[0];
10     p2 = gm1*(q2[3]-0.5f*ri*(q2[1]*q2[1]+q2[2]*q2[2]));
11     vol2 = ri*(q2[1]*dy - q2[2]*dx);
12     mu = 0.5f*((*adt1)+(*adt2))*eps;
13     f = 0.5f*(vol1* q1[0] + vol2* q2[0] ) + mu*(q1[0]-q2[0]);
14     res1[0] += f;
15     res2[0] -= f;
16     f = 0.5f*(vol1* q1[1] + p1*dy + vol2* q2[1] + p2*dy) + mu*(q1[1]-q2[1]);
17     res1[1] += f;
18     res2[1] -= f;
19     f = 0.5f*(vol1* q1[2] - p1*dx + vol2* q2[2] - p2*dx) + mu*(q1[2]-q2[2]);
20     res1[2] += f;
21     res2[2] -= f;
22     f = 0.5f*(vol1*(q1[3]+p1) + vol2*(q2[3]+p2) ) + mu*(q1[3]-q2[3]);
23     res1[3] += f;
24     res2[3] -= f;
25 }

```

Listing 1: Definition of the res.calc kernel with reals represented as single precision floating point numbers. Note the type signature. The kernel requires the element of the dataset x associated with each of the two nodes of the edge we are currently processing and the q , adt and res elements of the two cells associated with the current edge. Note that the res set is updated by incrementing. The important part of this are the data requirements of the kernel and not the exact meaning of the arithmetic operations. The variables $gm1$ and eps are global constants that do not need to be passed in explicitly.

2.2.2 Indirection maps

Having defined the data sets and the kernels, we now need to define the indirection maps that express the connectivity of the mesh and the relationships between the elements of the mesh. Airfoil has five such maps called: *edge*, *cell*, *ecell*, *bedge*, *becell*. They are presented in figure 2.4.

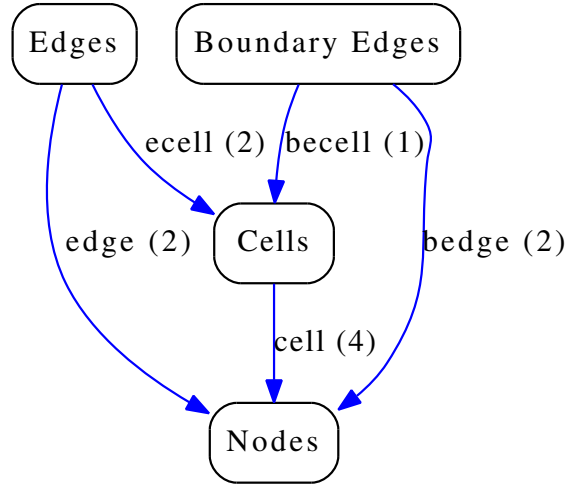


Figure 2.4: Diagram showing the maps between the mesh elements. The dimension of the map is shown in parentheses next to the name. Thus the map *edge* relating edges to nodes with dimension 2 means that for each edge, there are two nodes associated with it.

Having specified the data sets, indirection maps and kernels, the application of a kernel on an element is performed by looking up the mesh elements that element is associated with through the indirection maps and using those to access the data sets required by the kernel. For example, the invocation of the *res_calc* kernel defined in Listing 1 can be done with the following line of C:

```

res_calc(
    &x[2*edge[2*i]], &x[2*edge[2*i+1]], &q[4*ecell[2*i]],
    &q[4*ecell[2*i+1]], &adt[ecell[2*i]], &adt[ecell[2*i+1]],
    &res[4*ecell[2*i]], &res[4*ecell[2*i+1]]
);

```

Recall that *res_calc* operates on edges, and correlate the arguments to the type signature in Listing 1. There is a double level of indirection going on

here. i is the number of the edge we are currently processing (i ranges in $[0..number_of_edges - 1]$). As described in the section on mesh representation, the two nodes corresponding to the edge are stored at indices $2 * i$ and $2 * i + 1$ of the *edge* map. For each of those nodes, *res_calc* requires the corresponding element in the *x* data set. Recall from table 2.1 that the *x* set has a dimension of 2. Therefore the node numbers acquired from *edge* $[2 * i]$ and *edge* $[2 * i + 1]$ are multiplied by 2 and used as indices into the array *x* to access the correct data. Similarly for the rest of the arguments.

The complete iteration step in a sequential implementation of Airfoil is shown in Listing 2. The old values of q are stored in *q_old* and the inner loop runs twice before saving the solution again. The metric *rms* is computed in each iteration that is used to measure the convergence of the solution. The variables *ncell*, *nedge*, *nbedge* represent the number of cells, the number of edges and the number of boundary edges respectively.

A run of the sequential version in Listing 2 on a mesh with 721801 nodes, 1438600 edges, 2800 boundary edges and 720000 cells on an Intel Core2 Duo CPU at 2.8 GHz takes about 183 seconds to complete 1000 iterations. The time spent in each kernel is presented in table 2.3. It is evident that the computation is dominated by the *res_calc* and *adt_calc* kernels. In this project we will be concentrating on accelerating the *res_calc* kernel because it is the most computationally intensive kernel and because it has the most complex data access patterns that make it the interesting case to study. Finding a way to accelerate *res_calc* would pave the way for accelerating any similar kernel.

| Kernel Name | Time spent (seconds) | Percentage of total time (%) |
|-------------|----------------------|------------------------------|
| save_soln | 6.35 | 3.47 |
| adt_calc | 71.55 | 39.13 |
| res_calc | 81.57 | 44.62 |
| bres_calc | 0.43 | 0.24 |
| update | 22.93 | 12.54 |

Table 2.3: Table showing the time spent in each kernel during a run of a single-threaded sequential version of Airfoil. The total run time is 183 seconds.

```

1  int niter = 1000;
2  float rms = 0.0;
3  for(int iter=1; iter<=niter; iter++) {
4      for (int i = 0; i < ncell; ++i) {
5          save_soln(&q[4*i], &qold[4*i]);
6      }
7      for(int k=0; k<2; k++) {
8
9          for (int i = 0; i < ncell; ++i) {
10             adt_calc(&x[2*cell[4*i]], &x[2*cell[4*i+1]],
11                     &x[2*cell[4*i+2]], &x[2*cell[4*i+3]],
12                     &q[4*i], &adt[i]
13             );
14         }
15
16         for (int i = 0; i < nedge; ++i) {
17             res_calc(&x[2*edge[2*i]], &x[2*edge[2*i+1]],
18                     &q[4*ecell[2*i]], &q[4*ecell[2*i+1]],
19                     &adt[ecell[2*i]], &adt[ecell[2*i+1]],
20                     &res[4*ecell[2*i]], &res[4*ecell[2*i+1]]
21             );
22         }
23
24         for (int i = 0; i < nbedge; ++i) {
25             bres_calc(&x[2*bedge[2*i]], &x[2*bedge[2*i+1]],
26                      &q[4*becell[i]], &adt[becell[i]],
27                      &res[4*becell[i]], &bound[i]
28             );
29         }
30
31         rms = 0.0;
32         for (int i = 0; i < ncell; ++i) {
33             update(&qold[4*i], &q[4*i], &res[4*i], &adt[i], &rms);
34         }
35     }
36     rms = sqrt(rms/(float) ncell);
37     if (iter%100 == 0)
38         printf(" %d %10.5e \n",iter,rms);
39 }

```

Listing 2: The iteration structure of Airfoil.

2.3 Hardware platform, Maxeler toolchain and the streaming model of computation

The toolchain we use for implementing the FPGA accelerator is the one developed and maintained by Maxeler Technologies. It consists of the MAX3 cards that contain a Xilinx Virtex-6 chip [3] and up to 48GB of DDR3 DRAM. These cards can be programmed through MaxCompiler[4], which provides a Java-compatible object-oriented API to specify the dataflow graph. MaxCompiler will then schedule the graph, i.e. it will insert buffers that will introduce the appropriate delays in the design that will ensure the correct values will reach the appropriate stages in the pipeline at the correct clock cycle.

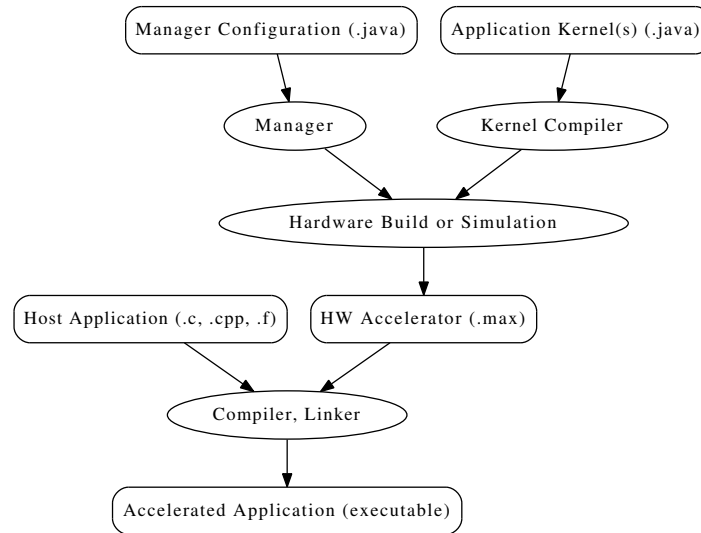


Figure 2.5: A diagram of the Maxeler toolchain. The data-flow graphs of the computational kernels are defined using a Java API. A manager connects multiple kernels together and handles the streaming to and from the kernels of data. These are combined by MaxCompiler and compiled into a .max file that can then be linked to a host C/C++ or Fortran application using standard tools (gcc, ld etc).

It will then produce a hardware design in VHDL that will then be further be compiled down to a binary bitstream that configures the FPGA by the Xilinx proprietary tools. The bitstream is then included in what is termed a *maxfile* that contains various other meta-data about the design

such as I/O stream names, named memory and register names, various run-time parameters etc. The maxfile can be linked against a normal C/C++ application using standard tools (gcc, ld etc). The interaction with the FPGA is performed by a low-level runtime: MaxCompilerRT and a driver layer: MaxelerOS. A diagram of the toolchain is shown in figure 2.5 [4].

Computational kernels in MaxCompiler have input streams that are pushed through a pipelined dataflow graph and some of them are output from the kernel. Programmatically, a hardware stream is seen as analogous to a variable in conventional programming languages. It's value potentially changes each cycle.

2.3.1 MaxCompiler example

We present a MaxCompiler design that computes a running 3-point average of a stream of floating point values (32 bits) in Listing 3.

```

1  public class MovingAverageKernel extends Kernel {
2
3      public MovingAverageKernel(KernelParameters parameters) {
4          super(parameters);
5          HWType flt = hwFloat(8,24);
6          HWVar x = io.input("x", flt );
7          HWVar x_prev = stream.offset(x, -1);
8          HWVar x_next = stream.offset(x, +1);
9          HWVar cnt = control.count.simpleCounter(32, N);
10         HWVar sel_nl = cnt > 0;
11         HWVar sel_nh = cnt < (N-1);
12         HWVar sel_m = sel_nl & sel_nh;
13         HWVar prev = sel_nl ? x_prev : 0;
14         HWVar next = sel_nh ? x_next : 0;
15         HWVar divisor = sel_m ? 3.0 : 2.0;
16         HWVar y = (prev+x+next)/divisor;
17         io.output("y" , y, flt);
18     }
19 }

```

Listing 3: A MaxCompiler definition of a kernel that computes a moving 3-point average with boundary conditions. Note that the arithmetic operators as well as the ternary if operator have been overloaded for HWVar objects that represent the value of a hardware stream.

MaxCompiler code is written in a Java-like language called MaxJ that provides overloaded operators such as $+$, $-$, $*$, $/$ and $? : .$ The example in Listing 3 creates a computational kernel that computes a stream of running 3-point averages, named y , from a stream of input values x . The HWVar class is the main representation of the value of a hardware stream at any clock cycle. HWVars always have a HWType that expresses the type of the stream (i.e. an integer, a floating point number, a 1-bit boolean value etc). The *stream.offset(x, -1)* and *stream.offset(x, +1)* expressions on lines 7 and 8 extract HWVars for the values of the stream on cycle in the past and one cycle in the future (note that this is internally done by creating implicit buffers, or FIFOs, and scheduling the pipelining accordingly). The ternary if operator $? :$ creates multiplexers in hardware that express choice. A Java API is provided that contains various useful design elements, such as counters (HWVars that increment their values in many configurable ways every cycle) that can be accessed through the control.count field.

The resulting dataflow graph can be seen in figure 2.6

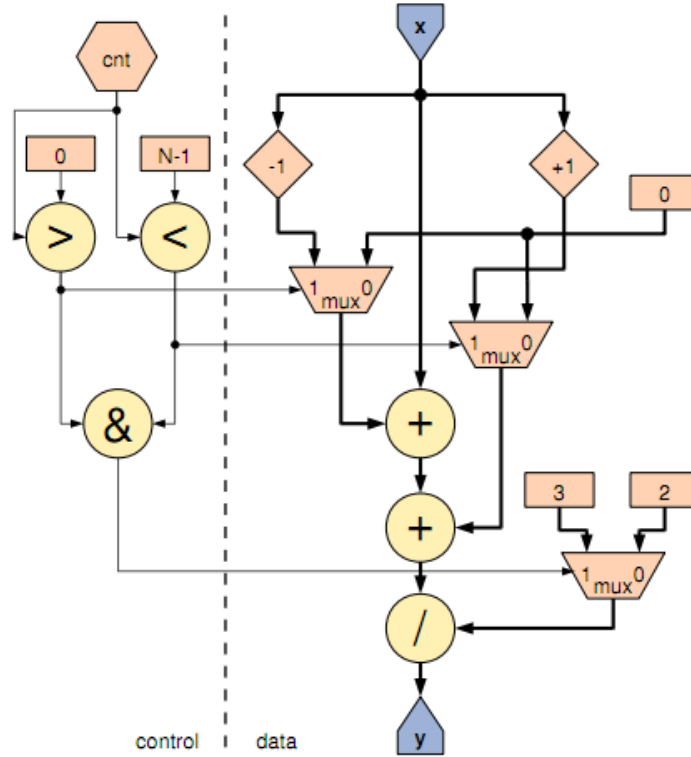


Figure 2.6: The dataflow graph resulting from the code in Listing 3

Kernel designs form part of a MaxCompiler design. The user also specifies a manager that describes the streaming connections between the kernels. A manager can be used to configure a design to stream data to and from the host through PCIe or from the DRAM that is attached to the FPGA. In the manager design, the user will instantiate the kernels and connect them up. Thus for the example in Listing 3 the manager might look like the one in Listing 4.

```

1  public class MovingAvgManager extends CustomManager {
2
3      public MovingAvgManager(MAXBoardModel board_model,
4                              boolean is_simulation, String name) {
5          super(is_simulation, board_model, name);
6          KernelBlock k
7              = addKernel(
8                  new MovingAverageKernel(makeKernelParameters("MovingAverageKernel"))
9              );
10
11         Stream x = addStreamFromHost("x");
12         k.getInput("x") <== x;
13
14         Stream y = addStreamToHost("y");
15         y <== k.getOutput("y");
16     }
17 }

```

Listing 4: Manager specification for a MovingAverageKernel that streams the input data "x" from the host and streams the output data "y" to the host. The <== operator means connect the right hand side stream to the left hand side stream. The above code instantiates the MovingAverageKernel, creates a stream called "x" from the host and connects it to the input stream "x" in the kernel. Then it creates a stream to the host called "y" and connects to it the output stream "y" from the kernel.

After we have specified a manger, we can build the design in order to create the .max file using the following lines of code:

```

public class MovingAvgHWBuilder {
    public static void main(String argv[]) {

        MovingAvgManager m
            = new MovingAvgManager(MAX3BoardModel.MAX3242A,
                                    false,
                                    "MovingAverage");

        m.build() ;
    }
}

```

This builds our design for a MAX3 card (containing a Xilinx Virtex6 FPGA)

using the "MovingAverage" name for the design.

Now that we have a .max file, we can interact with the FPGA from the host code by using the MaxCompilerRT API, an example of which is shown in Listing 5. In order to use the FPGA we must initialise the maxfile as in line 14 and open the device (line 15). The actual streaming to and from the FPGA is done using the max_run vararg function (line 22) where the arrays corresponding to the input data and the allocated space for the output data are specified. The MaxCompilerRT runtime and the MaxelerOS drivers handle the low-level details of PCIe streaming and interrupts.


```

1  #include<stdlib.h>
2  #include<stdint.h>
3  #include<MaxCompilerRT.h>
4  #define DATA_SIZE 1024
5
6  int main(int argc, char* argv[]) {
7      char* device_name = "/dev/maxeler0";
8      max_maxfile_t* maxfile;
9      max_device_handle_t* device;
10     float *data_in, *data_out;
11
12     maxfile = max_maxfile_init_MovingAverage();
13     device = max_open_device(maxfile, device_name);
14
15     data_in = (float*)malloc(DATA_SIZE * sizeof(float));
16     data_out = (float*)malloc(DATA_SIZE * sizeof(float));
17
18     for (int i = 0; i < DATA_SIZE; ++i) {
19         data_in[i] = i;
20     }
21
22     max_run(device,
23             max_input("x", data_in, DATA_SIZE * sizeof(float)),
24             max_output("y", data_out, DATA_SIZE * sizeof(float)),
25             max_runfor("MovingAverageKernel", DATA_SIZE),
26             max_end());
27
28
29     for (int i = 0; i < DATA_SIZE; ++i) {
30         printf("data_out[%d] = %f\n", i, data_out[i]);
31     }
32
33     max_close_device(device);
34     max_destroy(maxfile);
35     return 0;
36
37 }

```

Listing 5: A sample host code using the MaxCompilerRT API for the C language.

2.3.2 Hardware

The MAX3 card we use provides 48GB of DRAM that can be accessed with a maximum bandwidth of 38GB/s and a PCIe connection to the host machine that achieves a maximum bandwidth of 2GB/s in both directions. The Virtex6 FPGA by Xilinx used in the MAX3 card has about 4MB of fast on-board block RAM that should not be confused with the external DRAM. The host machine can communicate with the card through the PCIe bus using the MaxCompilerRT API. The external DRAM will be used to store the bulk of the mesh data and therefore achieving maximum utilisation of it is be one of the focal points of this project.

2.3.3 Mesh Partitioning and halos

In computing the optimal memory layout for our application, we have to partition large meshes into partitions that fit in the block RAM of the FPGA. We use a popular and widely available set of tools called METIS developed by George Karypis [5] that uses state of the art techniques to partition meshes, graphs, hypergraphs and matrices according to various parameters like size, edge/hyperedge cut, minimising certain metrics etc. It is a highly robust and efficient tool that we use through its C API. Since an iteration over a mesh element may require data associated with another element, partitions have a set of elements called a *halo region* that consists of all the elements (cells, nodes, edges) that maybe accessed from another partition. Consider figure 2.7. The partitioning is shown with the red line. The four partitions share cells (shown in purple) and edges (shown in red). This presents a difficulty when computing an edge that uses cells in the halo region of another partition, since we are storing a single partition at a time on the FPGA. The method used to deal with this issue is called the *halo exchange mechanism* and it opens up a large design space, with decisions usually dependent on the hardware and communication frameworks.

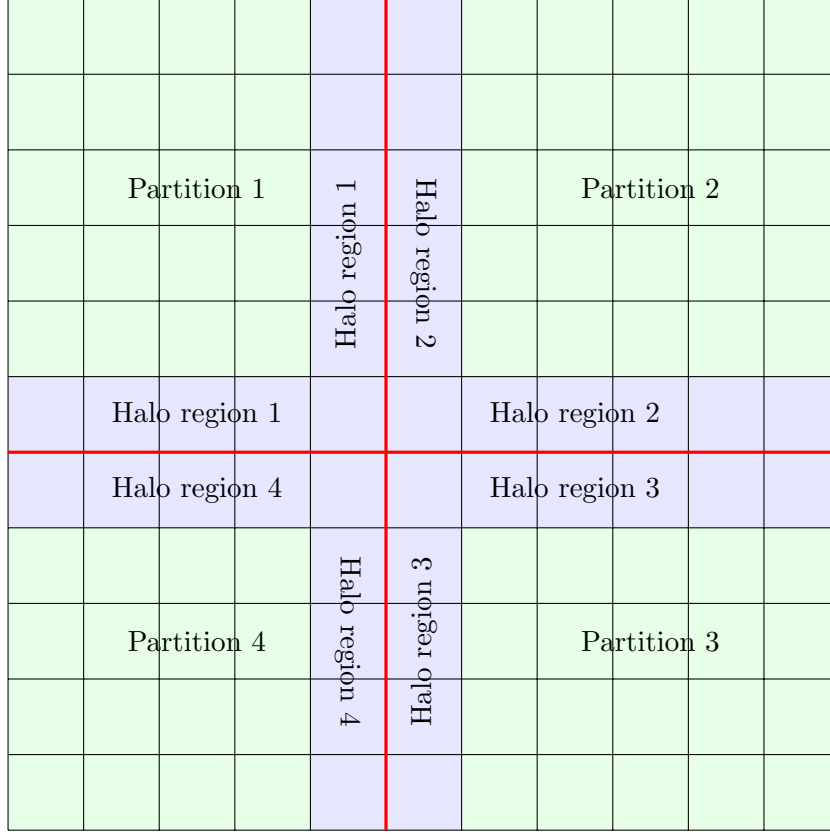


Figure 2.7: A mesh partitioned into 4 partitions, shown in green. The halo regions are shown in purple. Nodes, cells and edges belonging to the halo region can be accessed by another partition.

2.4 Previous work

Computation using unstructured meshes is widely used in many areas of engineering, not just in fluid dynamics, and there have been many attempts to augment the computation using accelerators. A recent trend has been to use the many cores available on Graphics Processing Units (GPUs) to launch thousands of threads in parallel, exploiting the parallel nature of many of these problems. Other hardware platforms include many-core architectures [1] and large parallel clusters [2]. More relevant to this project, there have been attempts to use FPGAs to accelerate such computations.

The principles of acceleration using FPGAs are quite different compared

to using GPUs or many-core systems. In the case of FPGA acceleration, the performance advantage comes from a custom, application specific, deeply pipelined datapath, often thousands of cycles deep that provides a throughput of one result per cycle. This argument for FPGA acceleration is widely accepted and is the source of the speedups achieved in all current attempts. Given this deep custom pipeline, it is a challenge for the developer to keep the pipeline fully utilised for the maximum amount of time and the techniques used to achieve this have been the main differentiating factors in the applications existing today. The most common approach has been to use on-chip block RAM memory to cache small parts of the mesh and operate on it, cache the results and write them back to main memory

M.T. Jones and K.Ramachandran [7] formulate the unstructured mesh computation as a sparse matrix problem, $Ax = y$ where A is a large sparse matrix representing the mesh and x is the vector that is being approximated. Their approach uses the conjugate gradient method to iteratively refine the approximation of the x vector. This involves, most importantly, a multiplication of the sparse matrix A with the vector x which forms the bulk of the computation and a subsequent refinement of the mesh and reconstruction of the sparse matrix. They formulate the problem as a sparse matrix-vector multiplication, whereas we are interested in iterating computational kernels over the mesh. Furthermore, they are concerned with mesh adaptation and the reconstruction of the sparse matrix in each iteration. We assume a static mesh specified at the beginning of the program.

Morishita et al. [8] examine the acceleration of CFD applications and in particular the use of on-chip block RAM resources to buffer the data in order to keep the arithmetic pipeline as full as possible. This is a more similar approach. However, their approach applies a constant stencil to a grid in 3D and tries to cache points in the grid that will be accessed in the next iteration, thus eliminating redundant accesses to the external memory. This caching/buffering is made possible by the fact that the stencil is of constant shape and thus the memory accesses can be predicted. In our application we have a 2D mesh that does not exhibit this property.

Sanchez-Roman et al. [9] present the acceleration of an airfoil-like unstructured mesh computation using FPGAs. Their solution uses two FPGAs on a single chip that perform different calculations and they identify the need to reason about computation and data access separately. They mention the need to partition larger meshes but they do not discuss techniques for partitioning or the issues arising from data dependencies across partitions. They mention the degradation in performance arising from the unstructured memory access patterns causing cache misses. We present a technique to

reorganise the mesh so as to facilitate more well-behaved memory accesses, allowing us to stream data to the datapath efficiently.

In another attempt, Sanchez-Roman et al. [10] recognise the update dependency that occurs during incrementing operation, similar to the ones in our `res_calc` kernel and work around it by adding an accumulator that correctly updates the required data sets. However, their design is used on comparatively small meshes of a maximum of 8000 nodes. Thus all the data can fit into the on-board memory of the FPGA, eliminating the need to consider partitioning issues. The applications we are concerned with usually have meshes of the order of 10^6 edges, which will definitely not fit on the on-chip block RAMs any time soon.

Chapter 3

Details and implementation

In this section we present our design, implementation and evaluation plan and provide the details of each stage.

3.1 Plan

We start by proposing various architectures for solving the problem and we propose a formal model for each one that allows us to predict the performance of the architecture. We then implement our scheme in order to provide real-world results and assess the feasibility of the implementation. This model will also allow us to pick the optimal values of the parameters of the application, thus maximizing performance and saving us the effort of using a trial and error approach to fine tune them. After that we proceed with the implementation of our chosen architecture.

During our work we realize that we have to partition the mesh into chunks that we can fit into the on-chip memory (called BRAM or block RAM) for processing. This requires us to think about and deal with data that overlap partitions (for example edges that begin in one partition and end in another). The set of shared data is known as the 'halo' of the partition and various halo exchange schemes exist. We use the ghost cell exchange mechanism, presented in [6].

After we have decided on the various parameters of the design (partition size, number of arithmetic pipelines, streaming responsibilities etc) we implement our design using MaxCompiler to produce a maxfile that we link to the Airfoil executable that we have modified to partition and layout the data in the decided way.

We can then compare our implementation against the theoretical model

developed earlier and track down and explain any and all discrepancies. Then we can compare our implementation against existing implementations that use Nvidia's OpenCL CUDA implementation.

Bibliography

- [1] MB Giles, GR Mudalige, Z Sharif, G Markall, PHJ Kelly,
Performance Analysis of the OP2 Framework on Many-core Architectures.
ACM SIGMETRICS Performance Evaluation Review, 38(4):9-15, March 2011
- [2] G.R Mudalige, MB Giles, C. Bertolli, P.H.J. Kelly,
Predictive Modeling and Analysis of OP2 on DistributedMemory GPU Clusters
PMBS '11 Proceedings of the second international workshop on Performance modeling, benchmarking and simulation of high performance computing systems Pages 3-4
- [3] Xilinx Inc.
Virtex-6 Family Overview
<http://www.xilinx.com/support/documentation/virtex-6.htm>
- [4] Maxeler Technologies
MaxCompiler White Paper
<http://www.maxeler.com/content/briefings/MaxelerWhitePaperMaxCompiler.pdf>
- [5] G. Karypis and V. Kumar.
A Fast and Highly Quality Multilevel Scheme for Partitioning Irregular Graphs
SIAM Journal on Scientific Computing, Vol. 20, No. 1, pp. 359392, 1999.
- [6] F. B. Kjolstad, M Snir
Ghost Cell Pattern
ParaPLoP '10 Proceedings of the 2010 Workshop on Parallel Programming Patterns
- [7] M. T. Jones, K. Ramachandran
Unstructured mesh computations on CCMs

Advances in Engineering Software - Special issue on large-scale analysis, design and intelligent synthesis environments Volume 31 Issue 8-9, Aug-Sept. 2000

- [8] H. Morishita, Y. Osana, N. Fujita, H. Amano
Exploiting memory hierarchy for a Computational Fluid Dynamics accelerator on FPGAs
ICECE Technology, 2008. FPT 2008. pp 193 - 200
- [9] Sanchez-Roman, D.; Sutter, G.; Lopez-Buedo, S.; Gonzalez, I.; Gomez-Arribas, F.J.; Aracil, J.; Palacios, F.;
High-Level Languages and Floating-Point Arithmetic for FPGABased CFD Simulations
Design & Test of Computers, IEEE, 2011, Volume: 28 Issue:4, pp 28 - 37
- [10] Sanchez-Roman, D.; Sutter, G.; Lopez-Buedo, S.; Gonzalez, I.; Gomez-Arribas, F.J.; Aracil, A.;
An Euler Solver Accelerator in FPGA for computational fluid dynamics applications
Proceedings of the 2011 VII Southern Conference on Programmable Logic Crdoba, Argentina April 13 - 15, 2011
- [11] Durbano, J.P.; Ortiz, F.E.;
FPGA-based acceleration of the 3D finite-difference time-domain method
12th Annual IEEE Symposium on Field-Programmable Custom Computing Machines, 2004. FCCM 2004.

Chapter 4

Appendix

4.1 Airfoil Kernel definitions in C

Even though we focused on accelerating the `res_calc` kernel, the other kernels are shown here for the sake of completeness.

```
void adt_calc(float *x1, float *x2, float *x3, float *x4, float *q, float *adt){
    float dx, dy, ri, u, v, c;
    ri = 1.0f/q[0];
    u = ri*q[1];
    v = ri*q[2];
    c = sqrt(gam*gm1*(ri*q[3]-0.5f*(u*u+v*v)));
    dx = x2[0] - x1[0];
    dy = x2[1] - x1[1];
    *adt = fabs(u*dy-v*dx) + c*sqrt(dx*dx+dy*dy);
    dx = x3[0] - x2[0];
    dy = x3[1] - x2[1];
    *adt += fabs(u*dy-v*dx) + c*sqrt(dx*dx+dy*dy);
    dx = x4[0] - x3[0];
    dy = x4[1] - x3[1];
    *adt += fabs(u*dy-v*dx) + c*sqrt(dx*dx+dy*dy);
    dx = x1[0] - x4[0];
    dy = x1[1] - x4[1];
    *adt += fabs(u*dy-v*dx) + c*sqrt(dx*dx+dy*dy);
    *adt = (*adt) / cfl;
}
```

```

void bres_calc(float *x1, float *x2, float *q1,
              float *adt1, float *res1, int *bound) {
    float dx, dy, mu, ri, p1, vol1, p2, vol2, f;
    dx = x1[0] - x2[0];
    dy = x1[1] - x2[1];
    ri = 1.0f/q1[0];
    p1 = gm1*(q1[3]-0.5f*ri*(q1[1]*q1[1]+q1[2]*q1[2]));
    if (*bound==1) {
        res1[1] += + p1*dy;
        res1[2] += - p1*dx;
    }
    else {
        vol1 = ri*(q1[1]*dy - q1[2]*dx);

        ri = 1.0f/qinf[0];
        p2 = gm1*(qinf[3]-0.5f*ri*(qinf[1]*qinf[1]+qinf[2]*qinf[2]));
        vol2 = ri*(qinf[1]*dy - qinf[2]*dx);

        mu = (*adt1)*eps;

        f = 0.5f*(vol1* q1[0] + vol2* qinf[0] ) + mu*(q1[0]-qinf[0]);
        res1[0] += f;
        f = 0.5f*(vol1* q1[1] + p1*dy + vol2* qinf[1] + p2*dy) + mu*(q1[1]-qinf[1]);
        res1[1] += f;
        f = 0.5f*(vol1* q1[2] - p1*dx + vol2* qinf[2] - p2*dx) + mu*(q1[2]-qinf[2]);
        res1[2] += f;
        f = 0.5f*(vol1*(q1[3]+p1) + vol2*(qinf[3]+p2) ) + mu*(q1[3]-qinf[3]);
        res1[3] += f;
    }
}

void res_calc(float *x1, float *x2, float *q1, float *q2,
             float *adt1, float *adt2, float *res1, float *res2) {

    float dx, dy, mu, ri, p1, vol1, p2, vol2, f;

    dx = x1[0] - x2[0];
    dy = x1[1] - x2[1];

    ri = 1.0f/q1[0];

```

```

    p1 = gm1*(q1[3]-0.5f*ri*(q1[1]*q1[1]+q1[2]*q1[2]));
    vol1 = ri*(q1[1]*dy - q1[2]*dx);

    ri = 1.0f/q2[0];
    p2 = gm1*(q2[3]-0.5f*ri*(q2[1]*q2[1]+q2[2]*q2[2]));
    vol2 = ri*(q2[1]*dy - q2[2]*dx);

    mu = 0.5f*((*adt1)+(*adt2))*eps;

    f = 0.5f*(vol1* q1[0] + vol2* q2[0] ) + mu*(q1[0]-q2[0]);
    res1[0] += f;
    res2[0] -= f;
    f = 0.5f*(vol1* q1[1] + p1*dy + vol2* q2[1] + p2*dy) + mu*(q1[1]-q2[1]);
    res1[1] += f;
    res2[1] -= f;
    f = 0.5f*(vol1* q1[2] - p1*dx + vol2* q2[2] - p2*dx) + mu*(q1[2]-q2[2]);
    res1[2] += f;
    res2[2] -= f;
    f = 0.5f*(vol1*(q1[3]+p1) + vol2*(q2[3]+p2) ) + mu*(q1[3]-q2[3]);
    res1[3] += f;
    res2[3] -= f;
}

void save_soln(float *q, float *qold){
    for (int n=0; n<4; n++) qold[n] = q[n];
}

void update(float *qold, float *q, float *res, float *adt, float *rms){
    float del, adti;

    adti = 1.0f/(*adt);

    for (int n=0; n<4; n++) {
        del = adti*res[n];
        q[n] = qold[n] - del;
        res[n] = 0.0f;
        *rms += del*del;
    }
}

```